Latency Distribution at 512 Tokens